

**AMENDMENTS TO THE CLAIMS:**

1. (Currently amended) A method of executing a linear algebra subroutine on a machine having at least one floating point unit (FPU) with one or more associated load/store units (LSU) to load data into and out of floating point registers (FRegs) of said FPU, said method comprising:

for an execution code controlling an operation of said floating point unit (FPU) performing a linear algebra subroutine execution involving three matrix operands in a level 3 nested loop matrix-matrix type kernel type operation (level 3 processing), inserting instructions to move moving data in a contiguous and stride one format into a cache providing data for said FPU for direct loading in a stride one manner into said FPU, so that said LSUs can load said data into said FRegs in an optimal manner before it is scheduled to be used in said linear algebra subroutine execution, said data being prefetched into said cache from a memory in a register block format, said register block format predetermined to reduce a number of provide a single data streams stream for each of said three operands involved in for a level 3 nested loop matrix-matrix type kernel type operation processing (e.g., said level 3 processing) to be three streams, thereby providing three streams of data for said level 3 processing, each said stream comprising contiguous data, and to allow a loading of these three streams into said FPU by said LSU,

said register block format comprising a data storage format wherein data is stored in blocks of size p-by-q, where p and q are small integers, meaning that p and q are sufficiently small so that ~~the~~ pieces of all data words of a block of size p-by-q of one or more of these blocks can be fitted into said FRegs as a result of one or more single instructions, said block comprising contiguous data to be moved stride one, and

wherein said three data streams comprise one stream of data of one matrix of said level 3 processing<sub>2</sub> as considered to be resident in said cache<sub>2</sub> and one stream each for data for two remaining matrix operands of said level 3 processing<sub>2</sub> as considered to respectively be residing in a memory or a cache having a level higher than or equal to a level of said cache.

2. (Previously presented) The method of claim 1, wherein said moving data is accomplished by scheduling move type instructions into time slots existing in a Level 3 Dense Linear Algebra Subroutine.

3. (Previously presented) The method of claim 1, wherein said linear algebra subroutine comprises a matrix multiplication operation.

4. (Previously presented) The method of claim 1, wherein said linear algebra subroutine comprises an equivalent of a subroutine from LAPACK (Linear Algebra PACKage).

5. (Previously presented) The method of claim 1, wherein said linear algebra subroutine invokes a BLAS Level 3 L1 cache kernel.

6. (Currently amended) An apparatus, comprising:

- a memory to store matrix data to be used for processing in a linear algebra program;
- a floating point unit (FPU) to perform said processing;
- a load/store unit (LSU) to load data to be processed by said FPU, said LSU loading said data into a plurality of floating point registers (FRegs); and
- a cache to store data from said memory and provide said data to said FRegs,

wherein said matrix data in said memory is moved by ~~having inserted moving~~  
~~instructions for said matrix data to be loaded~~ into said cache prior to a need for said data to be  
loaded by said LSU into said FRegs for said processing, said data being prefetched into said  
cache from said memory in a format predetermined to ~~reduce a number of~~ provide only three  
data streams for a level 3 nested loop matrix-matrix type kernel type operation (level 3 linear  
algebra processing) ~~to be three streams~~ and to allow a stride one loading of these streams into  
said FPU by said LSU (e.g., using SIMD (single instruction, multiple data)  $k > 1$ ,  $k$  being a  
number of data elements involved in said single instruction) ~~loading of these streams into said~~  
~~FPU by said LSU instructions,~~

wherein said format comprises a register block format wherein data is stored in blocks  
of size p-by-q, where p and q are small integers, meaning so that the pieces of all data words  
of a block of size p-by-q of one or more of these blocks can be fitted into said FRegs as a  
result of one or more single instructions, and

wherein said three data streams comprise one stream of data of one matrix of said  
level 3 linear algebra processing ~~is,~~ as considered to be resident in said cache, and one stream  
each for data for two remaining matrix operands of said level 3 linear algebra processing, as  
considered to respectively reside in a memory or a cache level higher than or equal to a level  
of said cache.

7. (Original) The apparatus of claim 6, wherein said linear algebra program comprises a  
matrix multiplication operation.

8. (Previously presented) The apparatus of claim 6, wherein said linear algebra program  
comprises an equivalent of a subroutine from LAPACK (Linear Algebra PACKage).

9. (Previously presented) The apparatus of claim 6, wherein said processing comprises invoking a BLAS Level 3 L1 cache kernel.

10. (Canceled)

11. (Previously presented) The apparatus of claim 6, wherein said moving instructions are inserted into time slots existing in a Level 3 Dense Linear Algebra Subroutine.

12. (Currently amended) A computer-readable storage medium tangibly embodying a program of machine-readable instructions executable by a digital processing apparatus to perform a method of executing linear algebra subroutines on a SIMD (single instruction, multiple data) machine having at least one floating point unit (FPU) with one or more associated load/store units (LSUs) to load data into and out of floating point registers (FRegs) of said FPU by way of a cache, said method comprising:

for an execution code controlling an operation of a floating point unit (FPU) performing a linear algebra subroutine execution, ~~inserting instructions to move~~ moving data into said cache providing said data into said FPU ~~before it was scheduled to be used for processing in said linear algebra subroutine,~~

wherein said data is prefetched into said cache from a memory in a format predetermined to ~~reduce a number of~~ provide three data streams for a level 3 nested loop matrix-matrix type kernel type operation (level 3 linear algebra processing) ~~to be three streams and to allow,~~ using a stride one (e.g., SIMD  $k > 1$  manner) loading of these three streams into said FPU by said LSUs with SIMD (single instruction, multiple data)  $k > 1$ ,  $k$  being a number of data elements involved in said single instruction) instructions,

wherein said format comprises a register block format wherein data is stored in blocks of size p-by-q, where p and q are small integers, meaning so that the pieces of all data words of a block of size p-by-q of one or more of these blocks can be fitted into said FRegs as a result of one or more single instructions, said block comprising contiguous data to be moved stride one, and

wherein said three data streams comprise one stream ~~as being of~~ data of one matrix of said level 3 linear algebra processing, as considered to be resident in said cache, and one stream each for data for two remaining streams, meaning one stream each for remaining two matrix operands of said level 3 linear algebra processing, as considered to respectively ~~that~~ reside in a memory or a cache having a level higher than or equal to a level of said cache.

13. (Previously presented) The computer-readable storage medium of claim 12, wherein said moving data is accomplished by inserting move type instructions into time slots existing in a Level 3 Dense Linear Algebra Subroutine.

14. (Previously presented) The computer-readable storage medium of claim 12, wherein said linear algebra subroutine comprises a matrix multiplication operation.

15. (Previously presented) The computer-readable storage medium of claim 12, wherein said linear algebra subroutine comprises an equivalent of a subroutine from LAPACK (Linear Algebra PACKage).

16. (Previously presented) The computer-readable storage medium of claim 12, wherein said linear algebra subroutine invokes a BLAS Level 3 L1 cache kernel.

Serial No. 10/671,889

Docket No. YOR920030170US1 (YOR.464)

17-20. (Canceled)